

On waiting time distributions associated with compound patterns in a sequence of multi-state trials

Kiyoshi Inoue · Sigeo Aki

Received: 5 April 2006 / Revised: 7 December 2006 / Published online: 18 August 2007
© The Institute of Statistical Mathematics, Tokyo 2007

Abstract In this article, waiting time distributions of compound patterns are considered in terms of the generating function of the numbers of occurrences of the compound patterns. Formulae for the evaluation of the generating functions of waiting time are given, which are very effective computational tools. We provide several viewpoints on waiting time problems associated with compound patterns and develop a general workable framework for the study of the corresponding distributions. The general theory is employed for the investigation of some examples in order to illustrate how the distributions of waiting time can be derived through our theoretical results.

Keywords Compound pattern · Scan · Run · Multi-state trials · Enumeration schemes · Conditional distribution · Probability function · Probability generating function · Double generating function

1 Introduction

Recently, there has been a great deal of interest in the development of the distribution theory of patterns in a sequence of multi-state trials (see [Antzoulakos 2001](#); [Inoue 2004](#); [Fu and Lou 2003](#); [Inoue and Aki 2002](#); [Fu and Chang 2002](#); [Hirano and Aki](#)

This research was partially supported by the ISM Cooperative Research Program (2006-ISM-CRP-2007).

K. Inoue (✉)
Faculty of Economics, Seikei University, 3-3-1 Kichijoji-Kitamachi,
Musasino-shi, Tokyo 180-8633, Japan
e-mail: kinoue@econ.seikei.ac.jp

S. Aki
Department of Mathematics, Faculty of Engineering, Kansai University, 3-3-35 Yamate-cho,
Suita-shi, Osaka 564-8680, Japan

2003; Stefanov 2000, 2003; Chryssaphinou and Papastavridis 1990; Han and Hirano 2003). Waiting time distributions of patterns are used effectively in a wide range of areas such as reliability, quality control and DNA sequence analysis (see Chao et al. 1995; Shmueli and Cohen 2000; Ewens and Grant 2001; Robin and Daudin 1999, 2001). However, it is very difficult to obtain the exact distributions, which usually involves hard probability theory and complicated mathematics. Even for the simple case where the underlying sequence is identically and independently distributed (i.i.d.) trials, many exact distributions remain unknown.

In this article, we study the waiting time distributions generated by compound patterns and investigate several aspects of the waiting time problems. The results presented here provide a proper framework for developing the exact distribution theory of compound patterns.

Let $\{Z_t, t \geq 1\}$ be a sequence of multi-state trials defined on the state space $\Gamma = \{0, 1, \dots, m\}$. According to Fu and Lou (2003) (see Fu and Chang 2002; Fu 1996), we will define a simple pattern and a compound pattern, respectively.

Definition 1 We say that ε is a simple pattern if ε is composed of a specified sequence of k states; i.e. $\varepsilon = (a_1, a_2, \dots, a_k)$, $a_i \in \Gamma$, $1 \leq i \leq k$ (the length of the pattern k is fixed, and the states in the pattern are allowed to be repeated). We identify ε with $\{\varepsilon\}$ and also call $\{\varepsilon\}$ the simple pattern. Let ε_1 and ε_2 be two simple patterns with size k_1 and k_2 respectively. We say that ε_1 and ε_2 are distinct if neither is a subsequence (segment) of the other.

Definition 2 We say that $c (\geq 2)$ simple patterns are distinct, if every two simple patterns among c simple patterns are distinct each other. We say that ε is a compound pattern if it is a union of $c (\geq 2)$ distinct simple patterns (a set of c distinct simple patterns). For the compound pattern $\varepsilon = \{\varepsilon_1, \dots, \varepsilon_c\}$, we define the occurrence of the compound pattern ε to be the occurrence of one of the simple patterns $\varepsilon_1, \dots, \varepsilon_c$.

Let $\varepsilon_i = \{\varepsilon_{i,j}, j = 1, \dots, c_i\}$, $i = 1, 2, \dots, \nu$, be compound patterns. We assume that all the simple patterns $\varepsilon_{i,j}$ ($i = 1, 2, \dots, \nu$, $j = 1, 2, \dots, c_i$) are distinct each other. For $i = 1, 2, \dots, \nu$, let $X_n^{\varepsilon_i}(\alpha_i)$ be the numbers of occurrences of compound pattern ε_i in Z_1, Z_2, \dots, Z_n under $\alpha_i (= N, O)$ counting, where the α_i represents the type of counting scheme employed; $\alpha_i = N$ will indicate the non-overlapping counting, $\alpha_i = O$ the overlapping counting. For $i = 1, 2, \dots, \nu$, we denote $E_{r_i}^{\varepsilon_i}(\alpha_i)$ by the event that the r_i compound patterns ε_i are observed in the sequence of multi-state trials under the α_i counting. Let $T_{\mathbf{r}}^{\varepsilon}(x; \alpha)$ be the waiting time for the x th occurrence of the event among $E_{r_i}^{\varepsilon_i}(\alpha_i)$ under the α_i counting, $i = 1, 2, \dots, \nu$, where $\varepsilon = (\varepsilon_1, \dots, \varepsilon_\nu)$, $\mathbf{r} = (r_1, \dots, r_\nu)$ and $\alpha = (\alpha_1, \dots, \alpha_\nu)$. Remark that each compound pattern ε_i is observed only r_i times, that is, after its r_i th occurrence we are no longer interested in the compound pattern ε_i and we are interested in when the remaining events occur. The random variable $T_{\mathbf{r}}^{\varepsilon}(1; \alpha)$ means the waiting time until at least one of the events $E_{r_i}^{\varepsilon_i}(\alpha_i)$, $i = 1, 2, \dots, \nu$ occurs. The random variable $T_{\mathbf{r}}^{\varepsilon}(2; \alpha)$ means the waiting time for the second occurrence among the events $E_{r_i}^{\varepsilon_i}(\alpha_i)$, $i = 1, 2, \dots, \nu$ occurs, where “the second occurrence” means the occurrence of another event excepting the first event among the events $E_{r_i}^{\varepsilon_i}(\alpha_i)$, $i = 1, 2, \dots, \nu$. Generally, the random variable $T_{\mathbf{r}}^{\varepsilon}(x; \alpha)$ means the waiting time for the x th occurrence among the events $E_{r_i}^{\varepsilon_i}(\alpha_i)$,

$i = 1, 2, \dots, v$ occurs. It is clear that $T_r^\varepsilon(1; \alpha) \leq T_r^\varepsilon(2; \alpha) \leq \dots \leq T_r^\varepsilon(v; \alpha)$. In the special case of $x = 1$ and $x = v$, the distributions of $T_r^\varepsilon(1; \alpha)$ and $T_r^\varepsilon(v; \alpha)$ are called *sooner waiting time distribution* and *later waiting time distribution*.

In Sect. 2, the distribution of the waiting time $T_r^\varepsilon(x; \alpha)$ is captured through the distribution of the sooner waiting time random variable $T_r^\varepsilon(1; \alpha)$. We investigate the relation between the distributions of waiting time $T_r^\varepsilon(x; \alpha)$ and the numbers $(X_n^{\varepsilon_1}(\alpha_1), \dots, X_n^{\varepsilon_v}(\alpha_v))$ of the occurrences of the compound patterns and proceed to derive formulae of the generating function of the waiting time $T_r^\varepsilon(x; \alpha)$ in terms of the generating function of $(X_n^{\varepsilon_1}(\alpha_1), \dots, X_n^{\varepsilon_v}(\alpha_v))$. Besides formulae of the generating function of the tail probability $P(T_r^\varepsilon(x; \alpha) > n)$ is established from the same viewpoint. Section 3 presents a discussion on conditional distribution of the waiting time $T_r^\varepsilon(x; \alpha)$, when the underlying sequence is a sequence of i.i.d. multi-state trials. In Sect. 4, the waiting time distributions for runs are explored and the generating functions are derived. As the special cases, the generalized birthday problem and the coupon collector’s problem are treated. Finally, Sect. 5 deals with the moving window scan statistics and linear/circular ratchet scan statistics for illustrative purposes.

2 General results

Let $\{Z_n, n \geq 1\}$ be a sequence of multi-state trials defined on the state space $\Gamma = \{0, 1, \dots, m\}$. Let $\varepsilon_i = \{\varepsilon_{i,j}, j = 1, \dots, c_i\}, i = 1, 2, \dots, v$, be compound patterns. As already mentioned in the introduction, we assume that all the simple patterns $\varepsilon_{i,j} (i = 1, 2, \dots, v, j = 1, 2, \dots, c_i)$ are distinct each other.

2.1 The waiting time for the x th occurrence of patterns

Let $T_r^\varepsilon(x; \alpha)$ be the waiting time for the x th occurrence of the event among $E_{r_i}^{\varepsilon_i}(\alpha_i), i = 1, 2, \dots, v$. The probability generating function and the double generating function of $T_r^\varepsilon(x; \alpha), r_i \geq 0, i = 1, 2, \dots, v$, will be denoted by $H_r^\varepsilon(t, x; \alpha)$ and $H^\varepsilon(t, z, x; \alpha)$, respectively, that is,

$$\begin{aligned}
 H_r^\varepsilon(t, x; \alpha) &= E \left[t^{T_r^\varepsilon(x; \alpha)} \right] = \sum_{n=0}^{\infty} P(T_r^\varepsilon(x; \alpha) = n) t^n, \\
 H^\varepsilon(t, z, x; \alpha) &= \sum_{r_1, \dots, r_v \geq 0} H_r^\varepsilon(t, x; \alpha) z_1^{r_1} \dots z_v^{r_v} \\
 &= \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} P(T_r^\varepsilon(x; \alpha) = n) t^n z_1^{r_1} \dots z_v^{r_v}.
 \end{aligned}$$

We will study the distribution of $T_r^\varepsilon(x; \alpha)$ through the distributions of the sooner waiting times $T_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(1; \alpha_{i_1}, \dots, \alpha_{i_j}), 1 \leq i_1 < \dots < i_j \leq v, j = v-x+1, \dots, v$. The key point for establishing the results is the relationship between the waiting time

$T_r^\varepsilon(x; \alpha)$ and the waiting times $T_{r_i}^{\varepsilon_i}(1; \alpha_i)$ until the r_i th occurrence of the compound pattern ε_i under α_i counting, $i = 1, 2, \dots, v$,

$$\{T_r^\varepsilon(x; \alpha) \geq n + 1\} \text{ if and only if}$$

$$\bigcup_{u=v-x+1}^v \bigcup_{\substack{1 \leq i_1 < \dots < i_u \leq v, \\ \{i_{u+1}, \dots, i_v\} \subset \{1, \dots, v\} \setminus \{i_1, \dots, i_u\} \\ 1 \leq i_{u+1} < \dots < i_v \leq v}} \left\{ T_{r_{i_1}}^{\varepsilon_{i_1}}(1; \alpha_{i_1}) \geq n + 1, \dots, T_{r_{i_u}}^{\varepsilon_{i_u}}(1; \alpha_{i_u}) \geq n + 1, \right. \\ \left. T_{r_{i_{u+1}}}^{\varepsilon_{i_{u+1}}}(1; \alpha_{i_{u+1}}) \leq n, \dots, T_{r_{i_v}}^{\varepsilon_{i_v}}(1; \alpha_{i_v}) \leq n \right\}. \tag{1}$$

Proposition 1 *We have the following relations:*

$$P(T_r^\varepsilon(x; \alpha) = n) = \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \\ \times P(T_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(1; \alpha_{i_1}, \dots, \alpha_{i_j}) = n), \tag{2}$$

$$H_r^\varepsilon(t, x; \alpha) = \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} H_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}} \\ \times (t, 1; \alpha_{i_1}, \dots, \alpha_{i_j}). \tag{3}$$

Proof From the relationship (1), we have

$$P(T_r^\varepsilon(x; \alpha) \geq n + 1) \\ = \sum_{u=v-x+1}^v \sum_{\substack{1 \leq i_1 < \dots < i_u \leq v, \\ \{i_{u+1}, \dots, i_v\} \subset \{1, \dots, v\} \setminus \{i_1, \dots, i_u\} \\ 1 \leq i_{u+1} < \dots < i_v \leq v}} P(T_{r_{i_1}}^{\varepsilon_{i_1}}(1; \alpha_{i_1}) \geq n + 1, \dots, \\ T_{r_{i_u}}^{\varepsilon_{i_u}}(1; \alpha_{i_u}) \geq n + 1, T_{r_{i_{u+1}}}^{\varepsilon_{i_{u+1}}}(1; \alpha_{i_{u+1}}) \leq n, \dots, T_{r_{i_v}}^{\varepsilon_{i_v}}(1; \alpha_{i_v}) \leq n) \\ = \sum_{u=v-x+1}^v \sum_{j=u}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-u} \binom{j}{u} \\ \times P(T_{r_{i_1}}^{\varepsilon_{i_1}}(1; \alpha_{i_1}) \geq n + 1, \dots, T_{r_{i_j}}^{\varepsilon_{i_j}}(1; \alpha_{i_j}) \geq n + 1).$$

Interchanging the order of the above summation and making use of the identity $\sum_{u=v-x+1}^j (-1)^{j-u} \binom{j}{u} = (-1)^{j-v+x-1} \binom{j-1}{v-x}$ (see for example [Feller 1968](#)), we have

$$P(T_r^\varepsilon(x; \alpha) \geq n + 1) \\ = \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \\ \times P(T_{r_{i_1}}^{\varepsilon_{i_1}}(1; \alpha_{i_1}) \geq n + 1, \dots, T_{r_{i_j}}^{\varepsilon_{i_j}}(1; \alpha_{i_j}) \geq n + 1)$$

$$\begin{aligned}
 &= \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \\
 &\quad \times P\left(T_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(1; \alpha_{i_1}, \dots, \alpha_{i_j}) \geq n+1\right).
 \end{aligned}$$

Observing that $P(T_{\mathbf{r}}^{\boldsymbol{\varepsilon}}(x; \boldsymbol{\alpha}) = n) = P(T_{\mathbf{r}}^{\boldsymbol{\varepsilon}}(x; \boldsymbol{\alpha}) \geq n) - P(T_{\mathbf{r}}^{\boldsymbol{\varepsilon}}(x; \boldsymbol{\alpha}) \geq n+1)$, we have the Eq. (2). The second conclusion (3) of the proposition is derived immediately by multiplying both sides of (2) by t^n and summing up for all $n \geq 0$. \square

Easily we see that the expected value $E[T_{\mathbf{r}}^{\boldsymbol{\varepsilon}}(x; \boldsymbol{\alpha})]$ can be captured through the expected values of the sooner waiting time $T_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(1; \alpha_{i_1}, \dots, \alpha_{i_j}), 1 \leq i_1 < \dots < i_j \leq v, j = v-x+1, \dots, v$;

$$\begin{aligned}
 E[T_{\mathbf{r}}^{\boldsymbol{\varepsilon}}(x; \boldsymbol{\alpha})] &= \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \\
 &\quad \times E[T_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(1; \alpha_{i_1}, \dots, \alpha_{i_j})].
 \end{aligned}$$

2.2 The relation between waiting times and numbers of occurrences of patterns

For $i = 1, 2, \dots, v$, let $X_n^{\varepsilon_i}(\alpha_i)$ be the numbers of occurrences of compound pattern ε_i in Z_1, Z_2, \dots, Z_n under $\alpha_i (= N, O)$ counting. Then, we define the probability generating function and the double generating function of $(X_n^{\varepsilon_1}(\alpha_1), \dots, X_n^{\varepsilon_v}(\alpha_v))$ by

$$\begin{aligned}
 \phi_n^{\boldsymbol{\varepsilon}}(\mathbf{z}; \boldsymbol{\alpha}) &= E\left[z_1^{X_n^{\varepsilon_1}(\alpha_1)} \dots z_v^{X_n^{\varepsilon_v}(\alpha_v)}\right] \\
 &= \sum_{x_1, \dots, x_v \geq 0} P(X_n^{\varepsilon_1}(\alpha_1) = x_1, \dots, X_n^{\varepsilon_v}(\alpha_v) = x_v) z_1^{x_1} \dots z_v^{x_v}, \\
 \Phi^{\boldsymbol{\varepsilon}}(\mathbf{z}, t; \boldsymbol{\alpha}) &= \sum_{n=0}^{\infty} \phi_n^{\boldsymbol{\varepsilon}}(\mathbf{z}; \boldsymbol{\alpha}) t^n \\
 &= \sum_{n=0}^{\infty} \sum_{x_1, \dots, x_v \geq 0} P(X_n^{\varepsilon_1}(\alpha_1) = x_1, \dots, X_n^{\varepsilon_v}(\alpha_v) = x_v) z_1^{x_1} \dots z_v^{x_v} t^n,
 \end{aligned}$$

respectively. Clearly, the probability generating function and double generating function of $(X_n^{\varepsilon_{i_1}}(\alpha_{i_1}), \dots, X_n^{\varepsilon_{i_j}}(\alpha_{i_j})), j = 1, 2, \dots, v$, can be expressed as

$$\begin{aligned}
 \phi_n^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}; \alpha_{i_1}, \dots, \alpha_{i_j}) &= \phi_n^{\boldsymbol{\varepsilon}}(\mathbf{z}; \boldsymbol{\alpha}) \Big|_{z_{iu}=1}, \quad \text{for } u \neq 1, 2, \dots, j, \\
 \Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j}) &= \Phi^{\boldsymbol{\varepsilon}}(\mathbf{z}, t; \boldsymbol{\alpha}) \Big|_{z_{iu}=1}, \quad \text{for } u \neq 1, 2, \dots, j.
 \end{aligned}$$

Let us elucidate the relation between the distributions of sooner waiting time $T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha})$ and $(X_n^{\mathcal{E}_1}(\alpha_1), \dots, X_n^{\mathcal{E}_v}(\alpha_v))$ in terms of the double generating functions. Notice that the dual relationship between the random variables $T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha})$ and $(X_n^{\mathcal{E}_1}(\alpha_1), \dots, X_n^{\mathcal{E}_v}(\alpha_v))$,

$$\{T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha}) > n\} \text{ if and only if } \{X_n^{\mathcal{E}_1}(\alpha_1) < r_1, \dots, X_n^{\mathcal{E}_v}(\alpha_v) < r_v\},$$

which gives the probability identity

$$P(T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha}) = n) = P(X_{n-1}^{\mathcal{E}_1}(\alpha_1) < r_1, \dots, X_{n-1}^{\mathcal{E}_v}(\alpha_v) < r_v) - P(X_n^{\mathcal{E}_1}(\alpha_1) < r_1, \dots, X_n^{\mathcal{E}_v}(\alpha_v) < r_v), \quad n, r_1, \dots, r_v \geq 1. \tag{4}$$

We set

$$P(T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha}) = 0) = \begin{cases} 1, & \text{if } r_i = 0 \text{ for some } i = 1, 2, \dots, v, \\ 0, & \text{otherwise.} \end{cases} \tag{5}$$

Lemma 1 *The double generating function $H^{\mathcal{E}}(t, \mathbf{z}, 1; \boldsymbol{\alpha})$ of $T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha})$ can be expressed in terms of double generating function $\Phi^{\mathcal{E}}(\mathbf{z}, t; \boldsymbol{\alpha})$.*

$$H^{\mathcal{E}}(t, \mathbf{z}, 1; \boldsymbol{\alpha}) = \frac{1}{\prod_{i=1}^v (1 - z_i)} \left(1 - \prod_{i=1}^v z_i (1 - t) \Phi^{\mathcal{E}}(\mathbf{z}, t; \boldsymbol{\alpha}) \right). \tag{6}$$

Proof By virtue of (4) and (5), we have

$$\begin{aligned} H^{\mathcal{E}}(t, \mathbf{z}, 1; \boldsymbol{\alpha}) &= \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} P(T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha}) = n) t^n z_1^{r_1} \cdots z_v^{r_v} \\ &= \sum_{r_1, \dots, r_v \geq 0} P(T_{\mathbf{r}}^{\mathcal{E}}(1; \boldsymbol{\alpha}) = 0) z_1^{r_1} \cdots z_v^{r_v} \\ &\quad + \sum_{r_1, \dots, r_v \geq 1} \sum_{n=1}^{\infty} \sum_{\substack{0 \leq i_j \leq r_j - 1 \\ j=1, \dots, v}} P(X_{n-1}^{\mathcal{E}_1}(\alpha_1) = i_1, \dots, X_{n-1}^{\mathcal{E}_v}(\alpha_v) = i_v) \\ &\quad \times t^n z_1^{r_1} \cdots z_v^{r_v} \\ &\quad - \sum_{r_1, \dots, r_v \geq 1} \sum_{n=1}^{\infty} \sum_{\substack{0 \leq i_j \leq r_j - 1 \\ j=1, \dots, v}} P(X_n^{\mathcal{E}_1}(\alpha_1) = i_1, \dots, X_n^{\mathcal{E}_v}(\alpha_v) = i_v) \\ &\quad \times t^n z_1^{r_1} \cdots z_v^{r_v}. \end{aligned}$$

Using the condition (5) and interchanging the orders of summation in the RHS, we get

$$\begin{aligned} &\sum_{r_1, \dots, r_v \geq 0} P(T_r^\mathcal{E}(1; \alpha) = 0) z_1^{r_1} \cdots z_v^{r_v} = \prod_{i=1}^v \frac{1}{1 - z_i} - \prod_{i=1}^v \frac{z_i}{1 - z_i}, \\ &\sum_{r_1, \dots, r_v \geq 1} \sum_{n=1}^{\infty} \sum_{\substack{0 \leq i_j \leq r_j - 1 \\ j=1, \dots, v}} P(X_{n-1}^{\varepsilon_1}(\alpha_1) = i_1, \dots, X_{n-1}^{\varepsilon_v}(\alpha_v) = i_v) t^n z_1^{r_1} \cdots z_v^{r_v} \\ &= \prod_{i=1}^v \frac{z_i}{1 - z_i} \sum_{n=1}^{\infty} \sum_{i_1, \dots, i_v \geq 0} P(X_{n-1}^{\varepsilon_1}(\alpha_1) = i_1, \dots, X_{n-1}^{\varepsilon_v}(\alpha_v) = i_v) t^n z_1^{i_1} \cdots z_v^{i_v} \\ &= \prod_{i=1}^v \frac{z_i}{1 - z_i} \sum_{n=1}^{\infty} \phi_{n-1}^\mathcal{E}(z; \alpha) t^n \quad \text{and} \\ &\sum_{r_1, \dots, r_v \geq 1} \sum_{n=1}^{\infty} \sum_{\substack{0 \leq i_j \leq r_j - 1 \\ j=1, \dots, v}} P(X_n^{\varepsilon_1}(\alpha_1) = i_1, \dots, X_n^{\varepsilon_v}(\alpha_v) = i_v) t^n z_1^{r_1} \cdots z_v^{r_v} \\ &= \prod_{i=1}^v \frac{z_i}{1 - z_i} \sum_{n=1}^{\infty} \phi_n^\mathcal{E}(z; \alpha) t^n. \end{aligned}$$

The proof is completed. □

Using the relation (6), we have the following theorem.

Theorem 1 *The double generating function $H^\mathcal{E}(t, z, x; \alpha)$ of $T_r^\mathcal{E}(x; \alpha)$ can be expressed in terms of the double generating functions $\Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j})$ of $(X_{r_{i_1}}^{\varepsilon_{i_1}}(\alpha_{i_1}), \dots, X_{r_{i_j}}^{\varepsilon_{i_j}}(\alpha_{i_j}))$, $1 \leq i_1 < \dots < i_j \leq v$, $j = v - x + 1, \dots, v$ as*

$$\begin{aligned} H^\mathcal{E}(t, z, x; \alpha) &= \frac{1}{\prod_{i=1}^v (1 - z_i)} \left(1 + \sum_{j=v-x+1}^v (-1)^{j-v+x} \binom{j-1}{v-x} \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j z_{i_u} \right. \\ &\quad \left. \times (1 - t) \Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j}) \right). \end{aligned}$$

Proof Substituting the Eq. (6) into the Eq. (3), we have

$$\begin{aligned} H^\mathcal{E}(t, z, x; \alpha) &= \sum_{r_1, \dots, r_v \geq 0} H_r^\mathcal{E}(t, x; \alpha) z_1^{r_1} \cdots z_v^{r_v} \\ &= \sum_{r_1, \dots, r_v \geq 0} \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \\ &\quad \times H_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(t, 1; \alpha_{i_1}, \dots, \alpha_{i_j}) z_1^{r_1} \cdots z_v^{r_v} \end{aligned}$$

$$\begin{aligned}
 &= \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \prod_{u \neq 1, \dots, j} \frac{1}{(1-z_{i_u})} \\
 &\quad \times H^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(t, z_{i_1}, \dots, z_{i_j}, 1; \alpha_{i_1}, \dots, \alpha_{i_j}) \\
 &= \frac{1}{\prod_{i=1}^v (1-z_i)} \sum_{j=v-x+1}^v \sum_{1 \leq i_1 < \dots < i_j \leq v} (-1)^{j-v+x-1} \binom{j-1}{v-x} \\
 &\quad \times \left(1 - \prod_{u=1}^j z_{i_u} (1-t) \Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j}) \right).
 \end{aligned}$$

Making use of the identity $\sum_{1 \leq i_1 < \dots < i_j \leq v} = \binom{v}{j}$ and using the relation $\sum_{j=v-x+1}^v \binom{v}{j} \binom{j-1}{v-x} (-1)^j = (-1)^{v+x-1}$ (see [Graham et al. 1994](#), p. 169), the proof is completed. \square

It should be noted that the double generating function $\Phi^{\mathbf{e}}(z, t, \alpha)$ is expressed in terms of the double generating functions of the sooner/later waiting time random variables;

$$\Phi^{\mathbf{e}}(z, t; \alpha) = \frac{1}{\prod_{i=1}^v z_i (1-t)} \left(1 - \prod_{i=1}^v (1-z_i) H^{\mathbf{e}}(t, z, 1; \alpha) \right), \tag{7}$$

$$\begin{aligned}
 \Phi^{\mathbf{e}}(z, t, \alpha) &= \frac{1}{\prod_{i=1}^v z_i (1-t)} \left(1 + \sum_{j=1}^v (-1)^j \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j (1-z_{i_u}) \right. \\
 &\quad \left. \times H^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(t, z_{i_1}, \dots, z_{i_j}, j; \alpha_{i_1}, \dots, \alpha_{i_j}) \right), \tag{8}
 \end{aligned}$$

where $H^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(t, z_{i_1}, \dots, z_{i_j}, j; \alpha_{i_1}, \dots, \alpha_{i_j})$ are the double generating functions of the later waiting time random variables $T_{r_{i_1}, \dots, r_{i_j}}^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(j; \alpha_{i_1}, \dots, \alpha_{i_j})$, $j = 1, 2, \dots, v$.

[Inoue and Aki \(2005b\)](#) have also given the inversion formulae (7) and (8) by a completely different technique. In the special case of $v = 1$, the results of the present subsection reduce to the ones derived by [Koutras \(1997\)](#). As by-products of the results presented in [Theorem 1](#), the generating function of the expected value $E[T_{\mathbf{r}}^{\mathbf{e}}(x; \alpha)]$ can be expressed as

$$\begin{aligned}
 &\sum_{r_1, \dots, r_v \geq 0} E[T_{\mathbf{r}}^{\mathbf{e}}(x; \alpha)] z_1^{r_1} \dots z_v^{r_v} \\
 &= \frac{1}{\prod_{i=1}^v (1-z_i)} \sum_{j=v-x+1}^v (-1)^{j-v+x-1} \binom{j-1}{v-x} \\
 &\quad \times \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j z_{i_u} \Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, 1; \alpha_{i_1}, \dots, \alpha_{i_j}).
 \end{aligned}$$

2.3 The tail probability of waiting time

The probability generating function and the double generating function of the tail probability $P(T_r^\epsilon(x; \alpha) > n)$ will be denoted by $\overline{H}_r^\epsilon(t, x; \alpha)$ and $\overline{H}^\epsilon(t, z, x; \alpha)$; i.e.,

$$\begin{aligned} \overline{H}_r^\epsilon(t, x; \alpha) &= \sum_{n=0}^{\infty} P(T_r^\epsilon(x; \alpha) > n) t^n, \\ \overline{H}^\epsilon(t, z, x; \alpha) &= \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} P(T_r^\epsilon(x; \alpha) > n) t^n z_1^{r_1} \dots z_v^{r_v}. \end{aligned}$$

Note that the above series are absolutely convergent at least $|t| < 1$ and $|z_i| < 1$, $i = 1, 2, \dots, v$. It is easy to see that the double generating function $\overline{H}^\epsilon(t, z, x; \alpha)$ can be captured through $\Phi^{\epsilon_{i_1}, \dots, \epsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j})$, $1 \leq i_1 < \dots < i_j \leq v$, $j = v - x + 1, \dots, v$. The next theorem provides the detail.

Theorem 2 *The double generating function $\overline{H}^\epsilon(t, z, x; \alpha)$ of the tail probability $P(T_r^\epsilon(x; \alpha) > n)$ can be expressed in terms of the double generating functions $\Phi^{\epsilon_{i_1}, \dots, \epsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j})$ of $(X_{r_{i_1}}^{\epsilon_{i_1}}(\alpha_{i_1}), \dots, X_{r_{i_j}}^{\epsilon_{i_j}}(\alpha_{i_j}))$, $1 \leq i_1 < \dots < i_j \leq v$, $j = v - x + 1, \dots, v$ as*

$$\begin{aligned} \overline{H}^\epsilon(t, z, x; \alpha) &= \frac{1}{\prod_{i=1}^v (1 - z_i)} \sum_{j=v-x+1}^v (-1)^{j-v+x-1} \binom{j-1}{v-x} \\ &\times \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j z_{i_u} \Phi^{\epsilon_{i_1}, \dots, \epsilon_{i_j}}(z_{i_1}, \dots, z_{i_j}, t; \alpha_{i_1}, \dots, \alpha_{i_j}). \end{aligned} \tag{9}$$

Proof We have

$$\begin{aligned} \overline{H}^\epsilon(t, z, x; \alpha) &= \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} P(T_r^\epsilon(x; \alpha) > n) t^n z_1^{r_1} \dots z_v^{r_v} \\ &= \sum_{r_1, \dots, r_v \geq 0} \frac{1 - H_r^\epsilon(t, x; \alpha)}{1 - t} z_1^{r_1} \dots z_v^{r_v} \\ &= \frac{1}{(1 - t) \prod_{i=1}^v (1 - z_i)} - \frac{1}{1 - t} H^\epsilon(t, z, x; \alpha). \end{aligned}$$

In view of Theorem 1, the proof is completed. □

It is noteworthy that the formula (9) produces expressions of $\Phi^\epsilon(z, t; \alpha)$ in terms of the double generating functions of the tail probabilities of sooner/later waiting time

random variables as

$$\begin{aligned} \Phi^{\mathbf{e}}(z, t; \boldsymbol{\alpha}) &= \prod_{i=1}^v \frac{(1 - z_i)}{z_i} \overline{H}^{\mathbf{e}}(t, z, 1; \boldsymbol{\alpha}), \\ \Phi^{\mathbf{e}}(z, t, \boldsymbol{\alpha}) &= \frac{1}{\prod_{i=1}^v z_i} \sum_{j=1}^v (-1)^{j-1} \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j (1 - z_{i_u}) \\ &\quad \times \overline{H}^{\mathbf{e}_{i_1, \dots, i_j}}(t, z_{i_1}, \dots, z_{i_j}, j; \alpha_{i_1}, \dots, \alpha_{i_j}), \end{aligned}$$

where $\overline{H}^{\mathbf{e}_{i_1, \dots, i_j}}(t, z_{i_1}, \dots, z_{i_j}, j; \alpha_{i_1}, \dots, \alpha_{i_j})$ are the double generating functions of the tail probabilities $P(T_{r_{i_1, \dots, i_j}}^{\mathbf{e}_{i_1, \dots, i_j}}(j; \alpha_{i_1}, \dots, \alpha_{i_j}) > n)$, $j = 1, 2, \dots, v$.

3 Conditional distributions

Let Z_1, Z_2, \dots, Z_n be a sequence of n i.i.d. random variables taking values in $\Gamma = \{0, 1, \dots, m\}$ and the probabilities $p_i = \Pr(Z_t = i)$, $1 \leq t \leq n$ and $i = 0, 1, \dots, m$. We are going to investigate the conditional distribution of the waiting time $T_{\mathbf{r}}^{\mathbf{e}}(x; \boldsymbol{\alpha})$, given the numbers $M_{n,i} = s_i$ ($0 \leq s_i \leq n$) of “ i ” ($i = 0, 1, \dots, m$) in the n i.i.d. trials. Since $M_{n,i}$ is a sufficient statistic for p_i ($i = 0, 1, \dots, m$), the conditional distribution which we are searching for does not depend on p_i ($i = 0, 1, \dots, m$). When the conditional distribution is considered and no confusion is likely to arise, we will use the notation $H_{\mathbf{r}}^{\mathbf{e}}(t, x, p_0, \dots, p_m; \boldsymbol{\alpha})$, $H^{\mathbf{e}}(t, z, x, p_0, \dots, p_m; \boldsymbol{\alpha})$, $\overline{H}_{\mathbf{r}}^{\mathbf{e}}(t, x, p_0, \dots, p_m; \boldsymbol{\alpha})$, $\overline{H}^{\mathbf{e}}(t, z, x, p_0, \dots, p_m; \boldsymbol{\alpha})$, $\Phi(z, t, p_0, \dots, p_m; \boldsymbol{\alpha})$ and $\phi_n(z, p_0, \dots, p_m; \boldsymbol{\alpha})$ instead of $H_{\mathbf{r}}^{\mathbf{e}}(t, x; \boldsymbol{\alpha})$, $H^{\mathbf{e}}(t, z, x; \boldsymbol{\alpha})$, $\overline{H}_{\mathbf{r}}^{\mathbf{e}}(t, x; \boldsymbol{\alpha})$, $\overline{H}^{\mathbf{e}}(t, z, x; \boldsymbol{\alpha})$, $\Phi(z, t; \boldsymbol{\alpha})$ and $\phi_n(z; \boldsymbol{\alpha})$. Again we write

$$H_{\mathbf{r}}^{\mathbf{e}}(t, x, p_0, \dots, p_m; \boldsymbol{\alpha}) = \sum_{n=0}^{\infty} P(T_{\mathbf{r}}^{\mathbf{e}}(x; \boldsymbol{\alpha}) = n) t^n, \tag{10}$$

$$H^{\mathbf{e}}(t, z, x, p_0, \dots, p_m; \boldsymbol{\alpha}) = \sum_{r_1, \dots, r_v \geq 0} H_{\mathbf{r}}^{\mathbf{e}}(t, x; \boldsymbol{\alpha}) z_1^{r_1} \dots z_v^{r_v}. \tag{11}$$

We will study the generating functions of the quantities

$$a_{\mathbf{r}}^{\mathbf{e}}(n, s, x; \boldsymbol{\alpha}) = \binom{n}{s_0, \dots, s_m} P(T_{\mathbf{r}}^{\mathbf{e}}(x; \boldsymbol{\alpha}) = n \mid M_{n,0} = s_0, \dots, M_{n,m} = s_m) \tag{12}$$

and

$$\bar{a}_{\mathbf{r}}^{\mathbf{e}}(n, s, x; \boldsymbol{\alpha}) = \binom{n}{s_0, \dots, s_m} P(T_{\mathbf{r}}^{\mathbf{e}}(x; \boldsymbol{\alpha}) > n \mid M_{n,0} = s_0, \dots, M_{n,m} = s_m), \tag{13}$$

where $s = (s_0, \dots, s_m)$.

Corollary 1 *The generating function of $a_{\mathbf{r}}^{\mathcal{E}}(n, s, n; \alpha)$ takes on the form*

$$\begin{aligned} & \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} \sum_{s_0 + \dots + s_m = n} a_{\mathbf{r}}^{\mathcal{E}}(n, s, x; \alpha) y_0^{s_0} \dots y_m^{s_m} t^n z_1^{r_1} \dots z_v^{r_v} \\ &= \frac{1}{\prod_{i=1}^v (1 - z_i)} \left(1 + \sum_{j=v-x+1}^v (-1)^{j-v+x} \binom{j-1}{v-x} \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j z_{i_u} (1-t) \right. \\ & \quad \left. \times \Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}} \left(z_{i_1}, \dots, z_{i_j}, \sum_{i=0}^m y_i t, \frac{y_0}{\sum_{i=0}^m y_i}, \dots, \frac{y_m}{\sum_{i=0}^m y_i}; \alpha_{i_1}, \dots, \alpha_{i_j} \right) \right). \end{aligned}$$

Proof Replacing $P(T_{\mathbf{r}}^{\mathcal{E}}(x; \alpha) = n)$ in (10) by

$$\begin{aligned} & P(T_{\mathbf{r}}^{\mathcal{E}}(x; \alpha) = n) \\ &= \sum_{s_0 + \dots + s_m = n} P(T_{\mathbf{r}}^{\mathcal{E}}(x; \alpha) = n \mid M_{n,0} = s_0, \dots, M_{n,m} = s_m) \\ & \quad \times P(M_{n,0} = s_0, \dots, M_{n,m} = s_m) \\ &= \sum_{s_0 + \dots + s_m = n} \binom{n}{s_0, \dots, s_m} p_0^{s_0} \dots p_m^{s_m} P(T_{\mathbf{r}}^{\mathcal{E}}(x; \alpha) = n \mid M_{n,0} = s_0, \dots, M_{n,m} = s_m) \end{aligned}$$

and exploiting the expression (11), we have

$$H_{\mathbf{r}}^{\mathcal{E}}(t, x, p_0, \dots, p_m; \alpha) = \sum_{n=0}^{\infty} \sum_{s_0 + \dots + s_m = n} a_{\mathbf{r}}^{\mathcal{E}}(n, s, x; \alpha) t^n p_0^{s_0} \dots p_m^{s_m}$$

or equivalently

$$\begin{aligned} H^{\mathcal{E}}(t, \mathbf{z}, x, p_0, \dots, p_m; \alpha) &= \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} \sum_{s_0 + \dots + s_m = n} a_{\mathbf{r}}^{\mathcal{E}}(n, s, x; \alpha) \\ & \quad \times p_0^{s_0} \dots p_m^{s_m} t^n z_1^{r_1} \dots z_v^{r_v}. \end{aligned}$$

Setting $p_i = y_i / \sum_{i=0}^m y_i$ ($i = 0, 1, \dots, m$) in the above expression, we get

$$\begin{aligned} & H^{\mathcal{E}}\left(t, \mathbf{z}, x, \frac{y_0}{\sum_{i=0}^m y_i}, \dots, \frac{y_m}{\sum_{i=0}^m y_i}; \alpha\right) \\ &= \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} \sum_{s_1 + \dots + s_v = n} a_{\mathbf{r}}^{\mathcal{E}}(n, s, x; \alpha) y_0^{s_0} \dots y_m^{s_m} \left(\frac{t}{\sum_{i=0}^m y_i}\right)^n z_1^{r_1} \dots z_v^{r_v}, \end{aligned}$$

which manifestly yields the desired result by replacing t by $\sum_{i=0}^m y_i t$. □

Similarly, in view of (13) we have the following corollary.

Corollary 2 *The generating function of $\bar{a}_r^\varepsilon(n, s, x; \alpha)$ takes on the form*

$$\begin{aligned} & \sum_{r_1, \dots, r_v \geq 0} \sum_{n=0}^{\infty} \sum_{s_0 + \dots + s_m = n} \bar{a}_r^\varepsilon(n, s, x; \alpha) y_0^{s_0} \dots y_m^{s_m} t^n z_1^{r_1} \dots z_v^{r_v} \\ &= \frac{1}{\prod_{i=1}^v (1 - z_i)} \sum_{j=v-x+1}^v (-1)^{j-v+x-1} \binom{j-1}{v-x} \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{u=1}^j z_{i_u} \\ & \times \Phi^{\varepsilon_{i_1}, \dots, \varepsilon_{i_j}} \left(z_{i_1}, \dots, z_{i_j}, \sum_{i=0}^m y_i t, \frac{y_0}{\sum_{i=0}^m y_i}, \dots, \frac{y_m}{\sum_{i=0}^m y_i}; \alpha_{i_1}, \dots, \alpha_{i_j} \right). \end{aligned}$$

4 Waiting time problems for runs

Let Z_1, Z_2, \dots, Z_n be a sequence of n i.i.d. random variables taking values in $\Gamma = \{0, 1, \dots, m\}$ with the probabilities $p_i = \Pr(Z_t = i)$, $1 \leq t \leq n$ and $i = 0, 1, \dots, m$. For $i = 1, 2, \dots, m$, let $\varepsilon_i = \underbrace{\{(i, i, \dots, i)\}}_{k_i}$ be the “ i ”-run of length k_i . In the litera-

ture, there are different ways of counting runs (see [Fu and Koutras 1994](#); [Balakrishnan and Koutras 2002](#)). It depends on the practical problem which way of counting should be adopted. The important and frequently used ways of counting runs are the “non-overlapping”, the “at least” and the “overlapping” scheme, which are called the Type I, II and III counting scheme, respectively (see [Balakrishnan and Koutras 2002](#); [Inoue and Aki 2005a](#)). As stated previously, the α_i represents the type of counting scheme employed for the “ i ”-run of length k_i ; $\alpha_i = N$ will indicate the non-overlapping counting, $\alpha_i = A$ the at least scheme and $\alpha_i = O$ overlapping one.

[Inoue and Aki \(2005a\)](#) derived the double generating function of $(X_n^{\varepsilon_1}(\alpha_1), \dots, X_n^{\varepsilon_m}(\alpha_m))$ as

$$\Phi^\varepsilon(z, t; \alpha) = \frac{1}{1 - p_0 t - \sum_{i=1}^m Q(z_i, p_i t; \alpha_i)}, \tag{14}$$

where

$$Q(z_i, p_i t; \alpha_i) = \begin{cases} \frac{p_i t - (p_i t)^{k_i} + (p_i t)^{k_i} z_i (1 - p_i t)}{1 - (p_i t)^{k_i}} & \alpha_i = N, \\ \frac{p_i t - (p_i t)^{k_i} (1 - z_i)}{1 - (p_i t)^{k_i} (1 - z_i)} & \alpha_i = A, \\ \frac{p_i t - (p_i t)^{k_i} (1 - z_i) - (p_i t)^2 z_i}{1 - p_i t z_i - (p_i t)^{k_i} (1 - z_i)} & \alpha_i = O, \end{cases} \tag{15}$$

for $i = 1, 2, \dots, m$.

Proposition 2 *The double generating function $H^\varepsilon(t, z, x; \alpha)$ of $T_r^\varepsilon(x; \alpha)$ is given by*

$$\begin{aligned}
 H^\varepsilon(t, z, x; \alpha) &= \frac{1}{\prod_{i=1}^m (1 - z_i)} \\
 &\times \left(1 + \sum_{j=m-x+1}^m (-1)^{j-m+x} \binom{j-1}{m-x} \sum_{1 \leq i_1 < \dots < i_j \leq m} \prod_{u=1}^j z_{i_u} (1-t) \right. \\
 &\quad \left. \times \frac{1}{1 - \left(1 - \sum_{u=1}^j p_{i_u}\right) t - \sum_{u=1}^j Q(z_{i_u}, p_{i_u} t; \alpha_{i_u})} \right). \tag{16}
 \end{aligned}$$

where $Q(z_i, p_i t, \alpha_i)$, $\alpha_i = N, A, O$, $i = 1, 2, \dots, m$ are as in (15).

Expanding the double generating function (16) in a multiple Taylor series around $z = \mathbf{0}$ and picking out the coefficient of $z_1 \cdots z_m$, we get the explicit expression for the probability generating function $H_{1, \dots, 1}^{\varepsilon_1, \dots, \varepsilon_m}(t, x; \alpha)$.

$$\begin{aligned}
 H_{1, \dots, 1}^{\varepsilon_1, \dots, \varepsilon_m}(t, x; \alpha) &= 1 + \sum_{j=m-x+1}^m (-1)^{j-m+x} \binom{j-1}{m-x} \\
 &\times \sum_{1 \leq i_1 < \dots < i_j \leq m} \frac{1-t}{1-t + \sum_{u=1}^j \frac{(p_{i_u} t)^{k_{i_u}} (1-p_{i_u} t)}{1-(p_{i_u} t)^{k_{i_u}}}},
 \end{aligned}$$

or equivalently

$$\begin{aligned}
 H_{1, \dots, 1}^{\varepsilon_1, \dots, \varepsilon_m}(t, x; \alpha) &= \sum_{j=m-x+1}^m (-1)^{j-m+x-1} \binom{j-1}{m-x} \\
 &\times \sum_{1 \leq i_1 < \dots < i_j \leq m} \frac{\sum_{u=1}^j \frac{(p_{i_u} t)^{k_{i_u}} (1-p_{i_u} t)}{1-(p_{i_u} t)^{k_{i_u}}}}{1-t + \sum_{u=1}^j \frac{(p_{i_u} t)^{k_{i_u}} (1-p_{i_u} t)}{1-(p_{i_u} t)^{k_{i_u}}}}.
 \end{aligned}$$

Needless to say, the probability generating function $H_{1, \dots, 1}^{\varepsilon_1, \dots, \varepsilon_m}(t, x; \alpha)$ is independent of α_i , $i = 1, 2, \dots, m$. For $x = 1$, $x = m$, the corresponding distributions are called *sooner/later geometric distributions of order (k_1, k_2, \dots, k_m)* , respectively. Further-

more, we obtain the explicit expression for the expected value of $T_{1,\dots,1}^{\varepsilon_1,\dots,\varepsilon_m}(x; \alpha)$ by differentiating $H_{1,\dots,1}^{\varepsilon_1,\dots,\varepsilon_m}(t, x; \alpha)$ with respect to t .

$$E \left[T_{1,\dots,1}^{\varepsilon_1,\dots,\varepsilon_m}(x; \alpha) \right] = \sum_{j=m-x+1}^v (-1)^{j-m+x-1} \binom{j-1}{v-x} \times \sum_{1 \leq i_1 < \dots < i_j \leq m} \frac{1}{\sum_{u=1}^j \frac{p_{i_u}^{k_{i_u}} (1 - p_{i_u})}{1 - p_{i_u}^{k_{i_u}}}}$$

Example 1 Birthday problems: Suppose that we interview people at random one by one, until we find r people with a common birthday. How many people should we have to interview? Specially, the case of $r = 2$ was investigated by many authors (see for example Johnson and Kotz (1977) and references therein). However, there are relatively few papers dealing with the general case ($r > 2$) and general arbitrary probabilities $p_1, \dots, p_{365}, p_1 + \dots + p_{365} = 1$. The results presented in this section will provide useful clues to the general birthday problems, since this problem can be captured through the distribution of sooner waiting time. The double generating function of the sooner waiting time $T_r^{\mathcal{E}}(1; \alpha)$ is expressed as

$$H^{\mathcal{E}}(t, z, 1; \alpha) = \frac{1}{\prod_{i=1}^{365} (1 - z_i)} \left(1 - \frac{z_1 \cdots z_{365} (1 - t)}{1 - (p_1 z_1 + \dots + p_{365} z_{365}) t} \right),$$

where $\varepsilon_i = \{(i)\}, \alpha_i = N, i = 1, 2, \dots, 365$.

Inoue and Aki (2005b) derive the probability generating function and the expected value of $T_r^{\mathcal{E}}(1; \alpha)$. Furthermore, in the special case where $p_i = 1/365, i = 1, 2, \dots, 365$, the expected value $E[T_{r,\dots,r}^{\varepsilon_1,\dots,\varepsilon_{365}}(1; \alpha)], r = 2, 3, \dots, 9$ is given numerically in the article.

Example 2 Coupon collector’s problems: Suppose that there are m distinct types of coupons bearing the numbers “1”, “2”, ..., “ m ” and that the coupon of type “ i ” is collected with probability $p_i (> 0), i = 1, 2, \dots, m, p_1 + \dots + p_m = 1$. We are interested in the total number of coupons until one collects x different types of coupons. When $x = m$, the later waiting time distribution is known as the coupon collector’s problem. In the special case where $p_i = 1/m, i = 1, 2, \dots, m$, the probability generating function and the expected value of $T_{1,\dots,1}^{\varepsilon_1,\dots,\varepsilon_m}(x; \alpha)$ are expressed as

$$H_{1,\dots,1}^{\varepsilon_1,\dots,\varepsilon_m}(t, x; \alpha) = \sum_{j=m-x+1}^m (-1)^{j-m+x-1} \binom{m}{j} \binom{j-1}{m-x} \frac{j t}{m - (m - j) t},$$

$$E[T_{1,\dots,1}^{\varepsilon_1,\dots,\varepsilon_m}(x; \alpha)] = \sum_{j=m-x+1}^m (-1)^{j-m+x-1} \binom{m}{j} \binom{j-1}{m-x} \frac{m}{j},$$

where $\varepsilon_i = \{(i)\}, \alpha_i = N, i = 1, 2, \dots, m$.

5 Scan statistics

We consider scan statistics, which are closely related to the sooner waiting time distribution of compound patterns. This section serves as an illustration of how the general theory presented in Sects. 2 and 3 can be employed for evaluating the distributions related to the scan statistics. Assume that the counting of all compound patterns ε_i ($i = 1, 2, \dots, \nu$) treated in this section are performed in the non-overlapping sense. We will suppress α_i ($i = 1, 2, \dots, \nu$) in the notations introduced in Sects. 2 and 3.

5.1 Moving window scan statistics

Let Z_1, Z_2, \dots, Z_n be a sequence of n i.i.d. random variables taking values in $\Gamma = \{0, 1, \dots, m\}$. Assume that $p_i = \Pr(Z_t = i)$, $t \geq 1$ and $i = 0, 1, \dots, m$. The scan statistic $S_n(w)$ of moving window of length w for the sequence Z_1, Z_2, \dots, Z_n is defined as

$$S_n(w) = \max \left\{ \sum_{j=i}^{i+w-1} Z_j : 1 \leq i \leq n - w + 1 \right\}.$$

We would like to study the unconditional probability $P(S_n(w) < r)$ and the conditional probability $P(S_n(w) < r | M_{n,0} = s_0, \dots, M_{n,m} = s_m)$ in terms of the distribution of the sooner waiting time. For illustrative purposes we consider the examples and proceed to the evaluation of the distributions by exploiting the results of Sects. 2 and 3.

Example 3 Assume that $w = 3, m = 2$ and $r = 5$. Then we can treat the unconditional probability $P(S_n(3) < 5)$ through the relation $P(S_n(3) < 5) = P(T_{1,1,1,1}^{\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4}(1) > n)$, where $\varepsilon_1 = \{(1, 2, 2)\}$, $\varepsilon_2 = \{(2, 1, 2)\}$, $\varepsilon_3 = \{(2, 2, 1)\}$, $\varepsilon_4 = \{(2, 2, 2)\}$. Easily we have $\Phi^{\varepsilon}(z, t) = P_1(z, t; \alpha) / P_0(z, t)$, where

$$\begin{aligned} P_1(z, t) &= 1 + p_2t + p_2(p_1 + p_2)t^2 - p_1p_2^2t^3 - p_1^2p_2^3t^5 \\ P_0(z, t) &= 1 - (p_0 + p_1)t - p_0p_2t^2 - p_2(p_1^2 + p_0p_2 + p_0p_1)t^3 + p_0p_1p_2^2t^4 \\ &\quad + p_0p_1^2p_2^3t^6 - p_1p_2^2t^3 \left(1 + p_1p_2t^2 \right) z_1 - p_1p_2^2t^3 z_2 \\ &\quad - (p_2t)^2 \left[1 - p_1t - p_1^2p_2t^3 \right] (p_1tz_3 + p_2tz_4). \end{aligned}$$

Using Theorem 2, the double generating function of the tail probability $P(T_r^{\varepsilon}(1) > n)$ can be expressed as

$$\bar{H}^{\varepsilon}(t, z, 1) = \prod_{i=1}^4 \frac{z_i}{(1 - z_i)} \frac{P_1(z, t)}{P_0(z, t)}. \tag{17}$$

Expanding (17) in a multiple Taylor series and picking out the coefficient of $t^n z_1 z_2 z_3 z_4$, we have

$$\sum_{n=0}^{\infty} P(S_n(3) < 5) t^n = \frac{1 + p_2 t + p_2(p_1 + p_2)t^2 - p_1 p_2^2 t^3 - p_1^2 p_2^3 t^5}{1 - (p_0 + p_1)t - p_0 p_2 t^2 - p_2(p_1^2 + p_0 p_1 + p_0 p_2)t^3 + p_0 p_1 p_2^2 t^4 + p_0 p_1^2 p_2^3 t^6}.$$

Example 4 (Continuation of Example 3) We consider the conditional probability $P(S_n(3) < 5 \mid M_{n,0} = s_0, M_{n,1} = s_1, M_{n,2} = s_2)$. Using Corollary 2 and the expression (17), we have

$$\sum_{n=0}^{\infty} \sum_{s_0+s_1+s_2=n} \bar{a}_{1,1,1,1}^{\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4}(n, \mathbf{s}) y_0^{s_0} y_1^{s_1} y_2^{s_2} t^n = \frac{1 + y_2 t + y_2(y_1 + y_2)t^2 - y_1 y_2^2 t^3 - y_1^2 y_2^3 t^5}{1 - (y_0 + y_1)t - y_0 y_2 t^2 - y_2(y_1^2 + y_0 y_1 + y_0 y_2)t^3 + y_0 y_1 y_2^2 t^4 + y_0 y_1^2 y_2^3 t^6}.$$

5.2 Linear and circular ratchet scan statistics

Let Z_1, Z_2, \dots, Z_n be a sequence of n i.i.d. random variables taking values in $\Gamma = \{1, 2, \dots, m\}$ with the probabilities

$$p_i = \Pr(Z_t = i), \quad 1 \leq t \leq n \text{ and } i = 1, 2, \dots, m.$$

For a given $w (< m)$, let $\varepsilon_i = \{(i), (i + 1), \dots, (i + w - 1)\}$, for $i = 1, 2, \dots, m - w + 1$ and $\varepsilon_i = \{(i), (i + 1), \dots, (m), (1), \dots, (w + i - 1 - m)\}$, for $i = m - w + 2, \dots, m$. Then we define $M_n(w) = \max_{1 \leq i \leq m-w+1} X_n^{\varepsilon_i}$ and define $M_n^c(w) = \max_{1 \leq i \leq m} X_n^{\varepsilon_i}$ (see Krauth 1999). The statistics $M_n(w)$ and $M_n^c(w)$ are called the linear/circular ratchet scan statistics, respectively. Specially, the statistic $M_n(1)$ is called disjoint statistic, when $\varepsilon_i = \{(i)\}$, $i = 1, \dots, m$. The linear/circular ratchet scan statistics are often applied to the problems in epidemiology (see Glaz et al. 2001). Krauth (1999) gives bounds for upper tail probabilities for the linear and circular ratchet scan statistics.

Let us consider the circular ratchet scan statistic. The probability $P(M_n^c(w) < r)$ is easily acquired by the distribution of the sooner waiting time $T_{r, \dots, r}^{\varepsilon_1, \dots, \varepsilon_m}(1)$. The identity $P(M_n^c(w) < r) = P(T_{r, \dots, r}^{\varepsilon_1, \dots, \varepsilon_m}(1) > n)$ will provide a way to compute the tail probability $P(M_n^c(w) < r)$. Observing that

$$\Phi^{\mathcal{E}}(z, t) = \frac{1}{1 - \sum_{i=1}^m p_i z_1^{I(i \in \varepsilon_1)} \dots z_m^{I(i \in \varepsilon_m)} t}$$

and using Theorem 2, the double generating function of the tail probability $P(T_r^{\mathcal{E}}(1) > n)$ can be expressed as

$$\begin{aligned} \overline{H}^{\mathcal{E}}(t, z, 1) &= \prod_{i=1}^v \frac{z_i}{(1-z_i)} \Phi^{\mathcal{E}}(z, t) \\ &= \prod_{i=1}^m \frac{z_i}{(1-z_i)} \frac{1}{1 - \sum_{i=1}^m p_i z_1^{I(i \in \varepsilon_1)} \cdots z_m^{I(i \in \varepsilon_m)} t}, \end{aligned} \quad (18)$$

where

$$I(v) = \begin{cases} 1, & v \text{ is true,} \\ 0, & \text{otherwise.} \end{cases}$$

Expanding (18) in a multiple Taylor series and picking out the coefficient of $t^n z_1^r \cdots z_m^r$, we can evaluate the tail probability $P(M_n^c(w) < r)$, which nowadays can be easily achieved by computer algebra systems.

Acknowledgements The authors wish to thank the editor and the referees for careful reading of our paper and helpful suggestions which led to improved results.

References

- Antzoulakos, D. L. (2001). Waiting times for patterns in a sequence of multistate trials. *Journal of Applied Probability*, 38, 508–518.
- Balakrishnan, N., Koutras, M. V. (2002). *Runs and scans with applications*. New York: Wiley.
- Chao, M. T., Fu, J. C., Koutras, M. V. (1995). Survey of reliability studies of consecutive- k -out-of- n : F & related systems. *IEEE Transactions on Reliability*, 40, 120–127.
- Chryssaphinou, O., Papastavridis, S. (1990). The occurrence of a sequence of patterns in repeated dependent experiments. *Theory of Probability and its Applications*, 35, 167–173.
- Ewens, W. J., Grant, G. R. (2001). *Statistical methods in bioinformatics: an introduction*. Heidelberg: Springer.
- Feller, W. (1968). *An introduction to probability theory and its applications*, Vol. I (3rd ed.). New York: Wiley.
- Fu, J. C. (1996). Distribution theory of runs and patterns associated with a sequence of multi-state trials. *Statistica Sinica*, 6, 957–974.
- Fu, J. C., Chang, Y. M. (2002). On probability generating functions for waiting time distributions of compound patterns in a sequence of multistate trials. *Journal of Applied Probability*, 39, 70–80.
- Fu, J. C., Koutras, M. V. (1994). Distribution theory of runs: a Markov chain approach. *Journal of the American Statistical Association*, 89, 1050–1058.
- Fu, J. C., Lou, W. Y. W. (2003). *Distribution theory of runs and patterns and its applications: a finite markov chain imbedding approach*. Singapore: World Scientific.
- Glaz, J., Naus, J., Wallenstein, S. (2001). *Scan statistics*. New York: Springer.
- Graham, R. L., Knuth, D. E., Patashnik, O. (1994). *Concrete mathematics* (2nd ed.). Reading: Addison-Wesley.
- Han, Q., Hirano, K. (2003). Sooner and later waiting time problems for patterns in Markov dependent trials. *Journal of Applied Probability*, 40, 73–86.
- Hirano, K., Aki, S. (2003). Number of occurrences of subpattern until the first appearance of a pattern and geometric distribution. *Statistics & Probability Letters*, 65, 259–262.
- Inoue, K. (2004). Joint distributions associated with patterns, successes and failures in a sequence of multi-state trials. *Annals of the Institute of Statistical Mathematics*, 56, 143–168.

- Inoue, K., Aki, S. (2002). Generalized waiting time problems associated with pattern in Polya's urn scheme. *Annals of the Institute of Statistical Mathematics*, 54, 681–688.
- Inoue, K., Aki, S. (2005a). A generalized Pólya urn model and related multivariate distributions. *Annals of the Institute of Statistical Mathematics*, 57, 49–59.
- Inoue, K., Aki, S. (2005b). On generating functions of waiting times and numbers of occurrences of compound patterns in a sequence of multi-state trials. Research Memorandum, No. 949, The Institute of Statistical Mathematics, Japan.
- Johnson, N. L., Kotz, S. (1977). *Urn models and their applications*. New York: Wiley.
- Koutras, M. V. (1997). Waiting times and number of appearances of events in a sequence of discrete random variables. In N. Balakrishnan (Ed.), *Advances in combinatorial methods and applications to probability and statistics* (pp. 363–384). Boston: Birkhauser.
- Krauth, J. (1999). Ratchet scan and disjoint scan statistics. In J. Glaz, & N. Balakrishnan (Eds.), *Scan statistics and applications* (pp. 67–96). Boston: Birkhauser.
- Robin, S., Daudin, J. J. (1999). Exact distribution of word occurrences in a random sequences of letters. *Journal of Applied Probability*, 36, 179–193.
- Robin, S., Daudin, J. J. (2001). Exact distribution of the distances between any occurrences of a set of words. *Annals of the Institute of Statistical Mathematics*, 53, 895–905.
- Shmueli, G., Cohen, A. (2000). Run-related probability functions applied to sampling inspection. *Technometrics*, 42, 188–202.
- Stefanov, V. T. (2000). On some waiting time problems. *Journal of Applied Probability*, 37, 756–764.
- Stefanov, V. T. (2003). The intersite distances between pattern occurrences in strings generated by general discrete—and continuous—time models: an algorithmic approach. *Journal of Applied Probability*, 40, 881–892.